

CPSC 330

Applied Machine Learning

Lecture 3: Machine Learning Fundamentals

UBC 2022-23

Instructor: Mathias Lécuyer

Imports

```
In [1]: 1 # import the libraries
2 import os
3 import sys
4
5 import graphviz
6 import IPython
7 import matplotlib.pyplot as plt
8 import numpy as np
9 import pandas as pd
10 from IPython.display import HTML
11 from sklearn.model_selection import train_test_split
12
13 sys.path.append("../code/.")
14 from plotting_functions import *
15
16 # Classifiers
17 from sklearn.tree import DecisionTreeClassifier, export_graphviz
18 from utils import *
19
20 %matplotlib inline
21
22 pd.set_option("display.max_colwidth", 200)
```

Learning outcomes

From this lecture, you will be able to:

- Understand model complexity and generalization
- Estimate the generalization error with data splitting

- Understand and use cross-validation
- Understand overfitting, underfitting, and related ML practice

Details:

Understand model complexity and generalization:

- explain how decision boundaries change with the `max_depth` hyperparameter;
- explain the concept of generalization;

Estimate the generalization error with data splitting

- appropriately split a dataset into train and test sets using `train_test_split` function;
- explain the difference between train, validation, test, and "deployment" data;
- identify the difference between training error, validation error, and test error;

Understand and use cross-validation

- explain cross-validation and use `cross_val_score` and `cross_validate` to calculate cross-validation error;

Understand overfitting, underfitting, and related ML practice

- recognize overfitting and/or underfitting by looking at train and test scores;
- explain why it is generally not possible to get a perfect test score (zero test error) on a supervised learning problem;
- describe the fundamental tradeoff between training score and the train-test gap;
- state the golden rule;
- start to build a standard recipe for supervised learning: train/test split, hyperparameter tuning with cross-validation, test on test set.

Announcements

- For assignments, setup, etc: use material from this year! Here: <https://github.com/UBC-CS/cpsc330-2022W2> (<https://github.com/UBC-CS/cpsc330-2022W2>). Links to older content are there if you want to explore, but not for deliverables.
- hw2 released. (Due next week Monday Jan 23 at 11:59pm.)
 - You are allowed to submit in pairs.
 - If you need to find a partner to work with, check out <https://piazza.com/class/lcg06c2ncl06el/post/5> (<https://piazza.com/class/lcg06c2ncl06el/post/5>).
- Advice on keeping up with the material
 - Practice! You will find some practice questions [here](https://ml-learn.mds.ubc.ca/) (<https://ml-learn.mds.ubc.ca/>).
 - Start early on homework assignments.
- Last day to withdraw without a W standing: January 23rd, 2022

Homework check in: <https://www.menti.com/als14k4vqbzx> (<https://www.menti.com/als14k4vqbzx>)

Recap

Last week, we introduced the following concepts:

- General idea of what Machine Learning is and its applications
- Basic terminology (such as features vs. target, unsupervised vs. supervised learning, classification vs. regression)
- Why is it a good idea to train a baseline classifier and how to do it (this included the basic steps of training using `scikit-learn`)
- Decision trees
- Parameters (e.g. decision tree rules) and hyperparameters (e.g., maximum tree depth)
- Decision boundary (today)

Pre-lecture videos

You were asked to watch 2 videos before coming to class, one on [generalization](https://youtu.be/iS2hsRRlc2M) (<https://youtu.be/iS2hsRRlc2M>), one on [data splitting](https://youtu.be/h2AEobwcUQw) (<https://youtu.be/h2AEobwcUQw>).

We will revisit the first one to discuss dicision boundaries, and quickly mention the second.

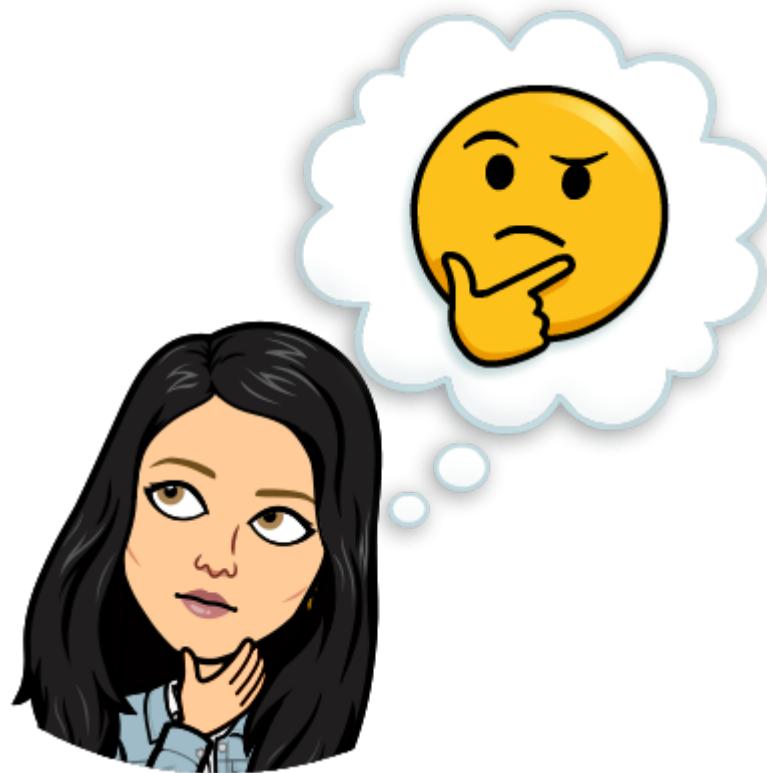
Generalization [[video \(https://youtu.be/iS2hsRRlc2M\)](https://youtu.be/iS2hsRRlc2M)]

Big picture and motivation

In machine learning, we want to glean information from labeled data so that we can label **new unlabeled** data.

For example, suppose we want to build a spam filtering system. We will take a large number of spam/non-spam messages from the past, learn patterns associated with spam/non-spam from them, and predict whether **a new incoming message** in someone's inbox is spam or non-spam based on these patterns.

So we want to learn from the past but ultimately we want to apply it on the future email messages.



How can we generalize from what we've seen to what we haven't seen?

In this lecture, we'll see how machine learning tackles this question.

Model complexity and training error

Let's examine the **decision boundary** of tree classifiers, and how it changes for different tree depths, to visualize what sort of examples will be classified as positive and negative.

In [2]:

```
1 # Toy quiz2 grade data
2 classification_df = pd.read_csv("../data/quiz2-grade-toy-classification")
3 classification_df.head(10)
```

Out[2]:

	ml_experience	class_attendance	lab1	lab2	lab3	lab4	quiz1	quiz2	
0	1		1	92	93	84	91	92	A+
1	1		0	94	90	80	83	91	not A+
2	0		0	78	85	83	80	80	not A+
3	0		1	91	94	92	91	89	A+
4	0		1	77	83	90	92	85	A+
5	1		0	70	73	68	74	71	not A+
6	1		0	80	88	89	88	91	A+
7	0		1	95	93	69	79	75	not A+
8	0		0	97	90	94	99	80	not A+
9	1		1	95	95	94	94	85	not A+

In [3]:

```
1 X = classification_df.drop(["quiz2"], axis=1)
2 y = classification_df["quiz2"]
3
4 X_subset = X[["lab4", "quiz1"]] # Let's consider a subset of the data
5 X_subset.head()
```

Out[3]:

	lab4	quiz1
0	91	92
1	83	91
2	80	80
3	91	89
4	92	85

In []:

```
1 X_subset = X[["lab4", "quiz1"]] # Let's consider a subset of the data
2 X_subset.head()
```

In the following model (decision stump), this decision boundary is created by asking one question.

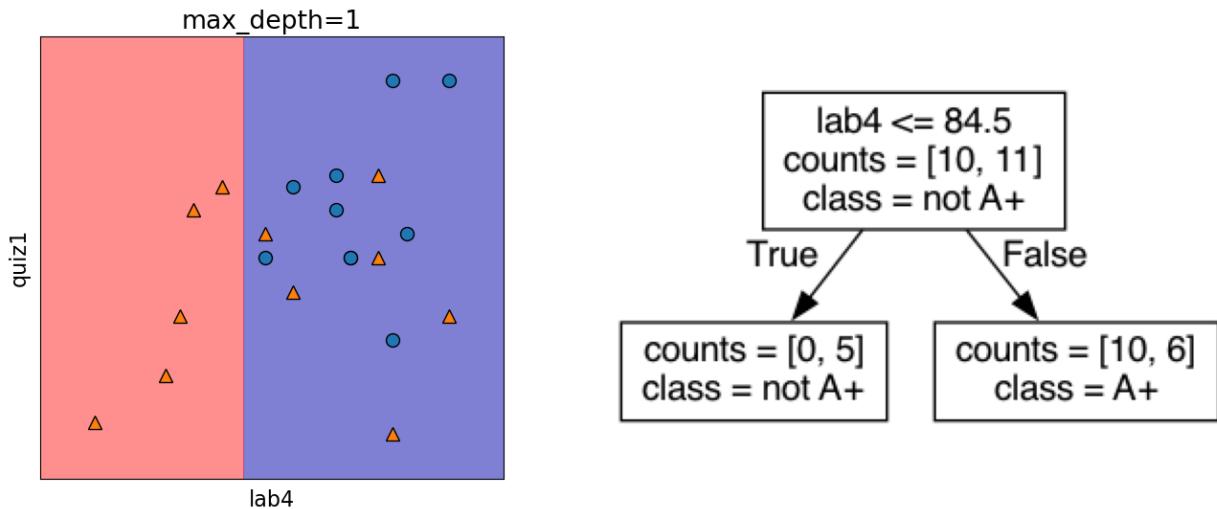
In [4]:

```

1 depth = 1
2 model = DecisionTreeClassifier(max_depth=depth)
3 model.fit(X_subset.to_numpy(), y)
4 model.score(X_subset.to_numpy(), y)
5 print("Error: %0.3f" % (1 - model.score(X_subset.to_numpy(), y)))
6 plot_tree_decision_boundary_and_tree(model, X_subset, y, x_label="lab4")

```

Error: 0.286



In []:

1

In the following model, this decision boundary is created by asking two questions.

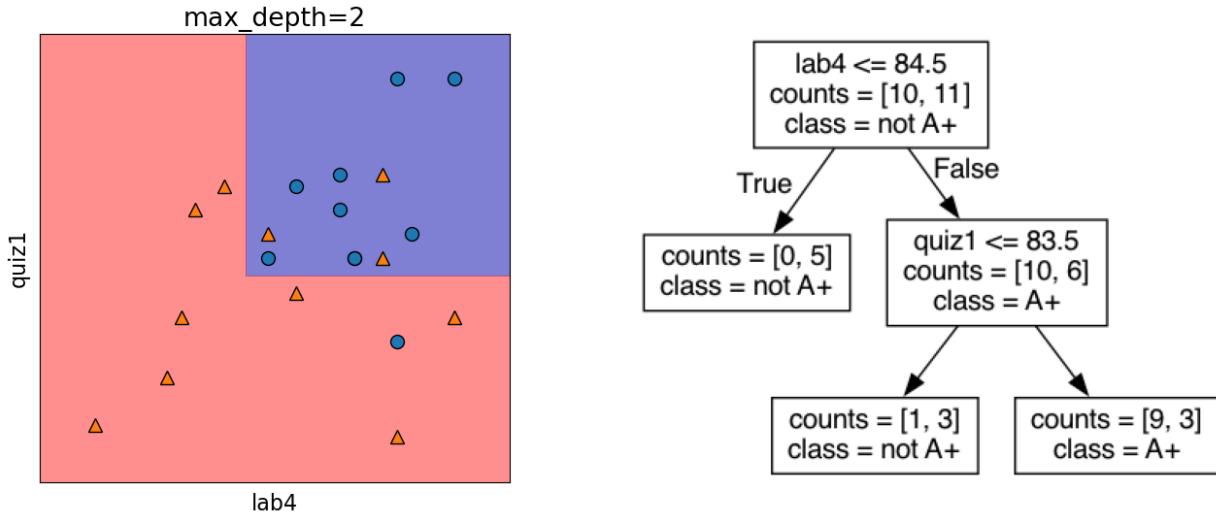
In [5]:

```

1 depth = 2
2 model = DecisionTreeClassifier(max_depth=depth)
3 model.fit(X_subset.to_numpy(), y)
4 model.score(X_subset.to_numpy(), y)
5 print("Error: %0.3f" % (1 - model.score(X_subset.to_numpy(), y)))
6 plot_tree_decision_boundary_and_tree(
    model, X_subset, y, x_label="lab4", y_label="quiz1"
8 )

```

Error: 0.190



Let's look at the decision boundary with depth = 4.

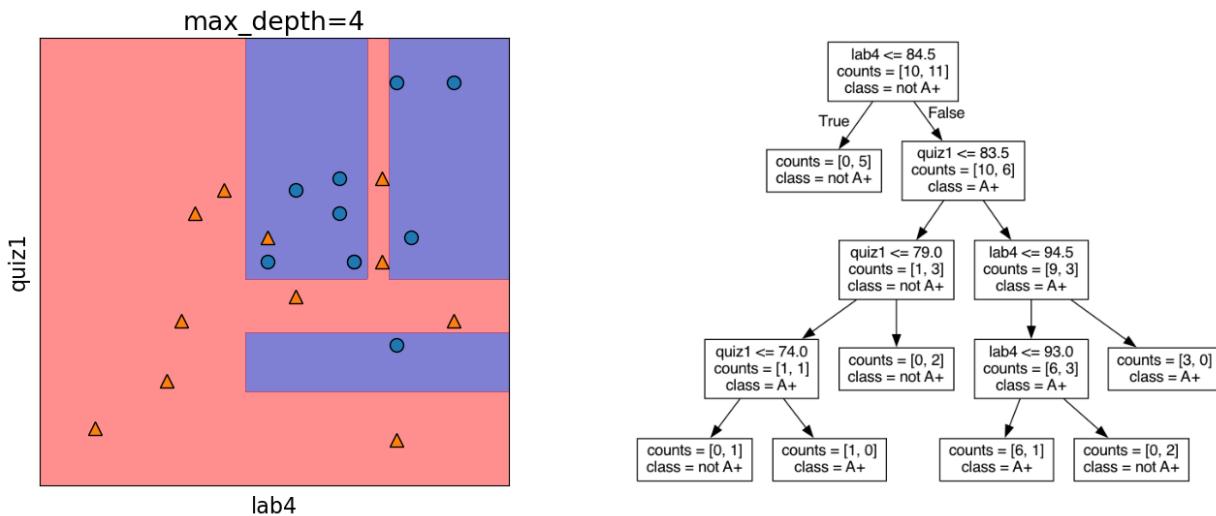
In [6]:

```

1 depth = 4
2 model = DecisionTreeClassifier(max_depth=depth)
3 model.fit(X_subset.to_numpy(), y)
4 model.score(X_subset.to_numpy(), y)
5 print("Error: %0.3f" % (1 - model.score(X_subset.to_numpy(), y)))
6 plot_tree_decision_boundary_and_tree(
    model, X_subset, y, x_label="lab4", y_label="quiz1"
8 )

```

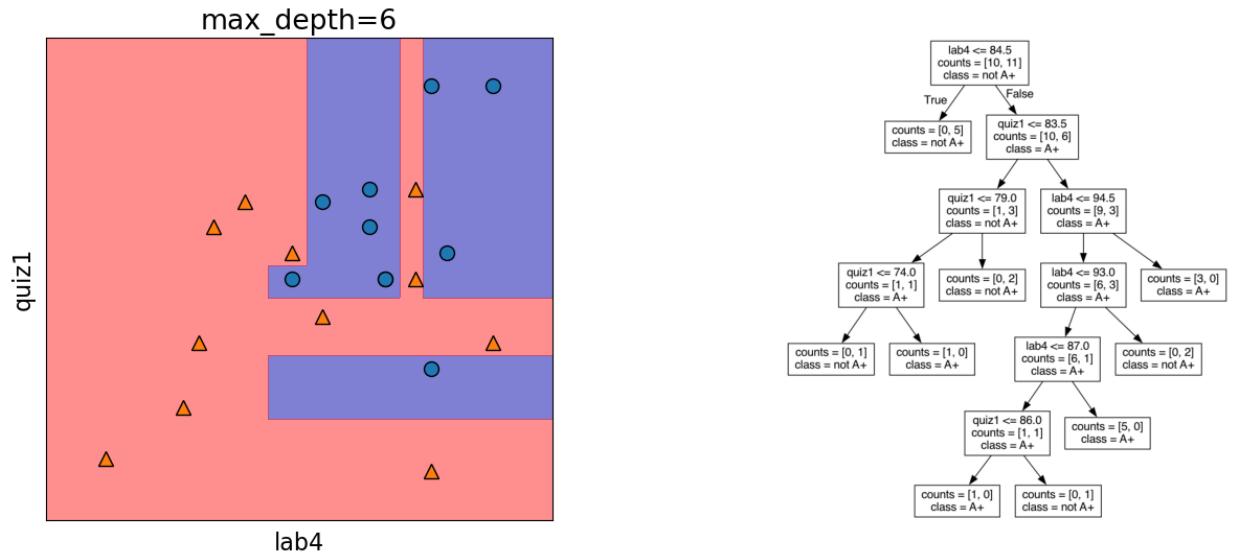
Error: 0.048



Let's look at the decision boundary with depth = 6.

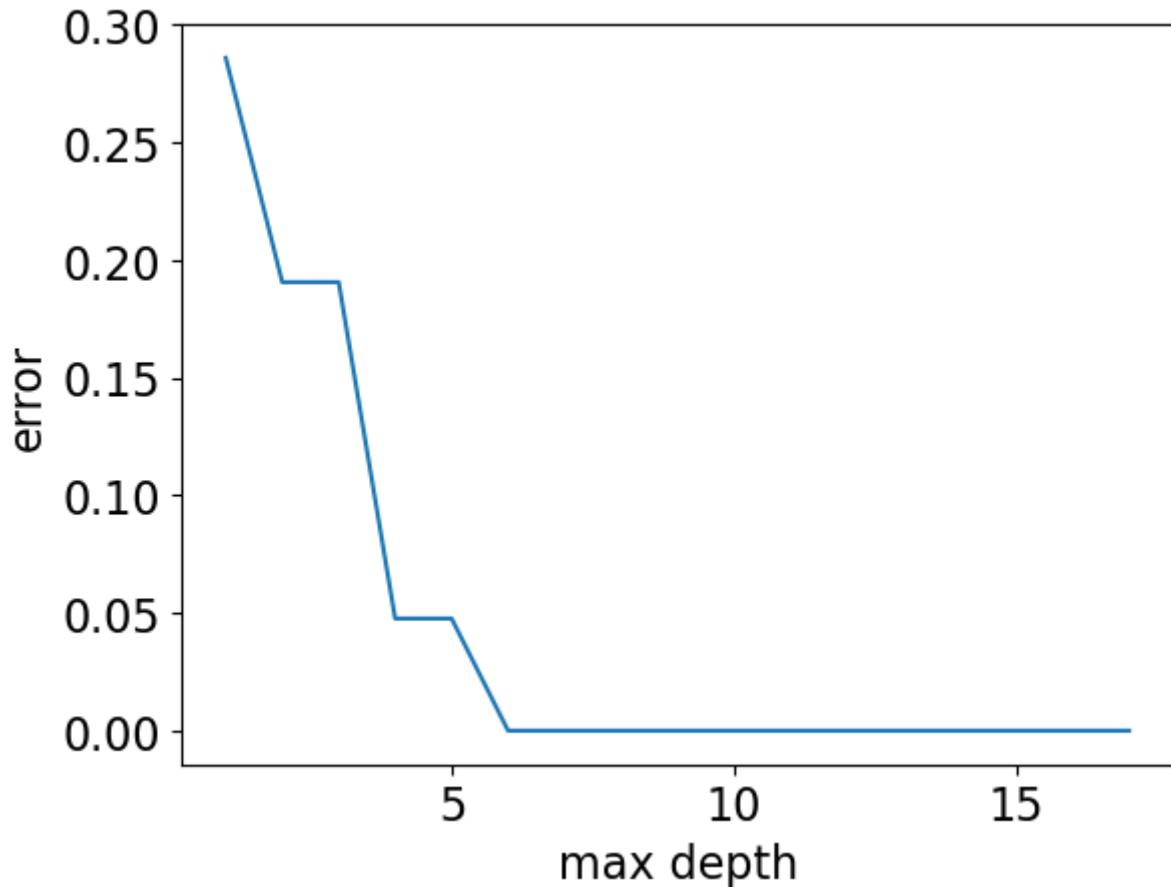
```
In [7]: 1 depth = 6
2 model = DecisionTreeClassifier(max_depth=depth)
3 model.fit(X_subset.to_numpy(), y)
4 model.score(X_subset.to_numpy(), y)
5 print("Error:  %0.3f" % (1 - model.score(X_subset.to_numpy(), y)))
6 plot_tree_decision_boundary_and_tree(
7     model, X_subset, y, x_label="lab4", y_label="quiz1"
8 )
```

Error: 0.000



In [8]:

```
1 max_depths = np.arange(1, 18)
2 errors = []
3 for max_depth in max_depths:
4     error = 1 - DecisionTreeClassifier(max_depth=max_depth).fit(X_subset,
5         X_subset, y
6     )
7     errors.append(error)
8 plt.plot(max_depths, errors)
9 plt.xlabel("max depth")
10 plt.ylabel("error");
```



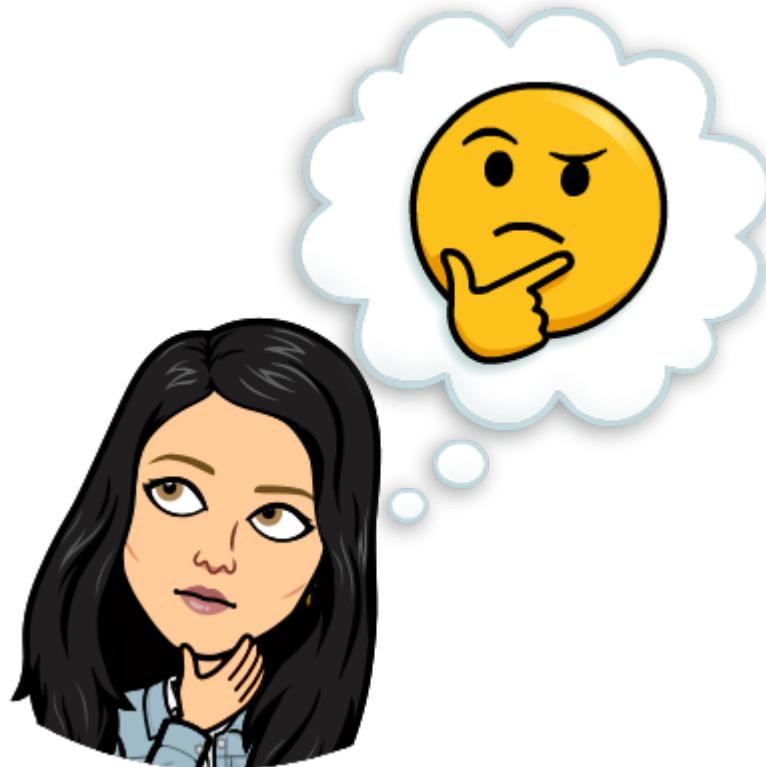
- Our model has 0% error for depths ≥ 6 !!
- But it's also becoming more and more specific and sensitive to the training data.
- Is it good or bad?

Although the plot above (complexity hyperparameter vs error) is more popular, we could also look at the same plot flip the y-axis, i.e., consider accuracy instead of error.

In []:

```
1 max_depths = np.arange(1, 18)
2 accuracies = []
3 for max_depth in max_depths:
4     accuracy = (
5         DecisionTreeClassifier(max_depth=max_depth).fit(X_subset, y).sc
6     )
7     accuracies.append(accuracy)
8 plt.plot(max_depths, accuracies)
9 plt.xlabel("max depth")
10 plt.ylabel("accuracy");
```

🤔 Eva's questions



At this point Eva is wondering about the following questions.

- How to pick the best depth?
- How can we make sure that the model we have built would do reasonably well on new data in the wild when it's deployed?
- Which of the following rules learned by the decision tree algorithm are likely to generalize better to new data?

Rule 1: If class_attendance == 1 then grade is A+.

Rule 2: If lab3 > 83.5 and quiz1 <= 83.5 and lab2 <= 88 then quiz2 grade is A+

Generalization: Fundamental goal of ML

To generalize beyond what we see in the training examples

We only have access to limited amount of training data and we want to learn a mapping function which would predict targets reasonably well for examples beyond this training data.

- Example: Imagine that a learner sees the following images and corresponding labels.

Generalizing to unseen data

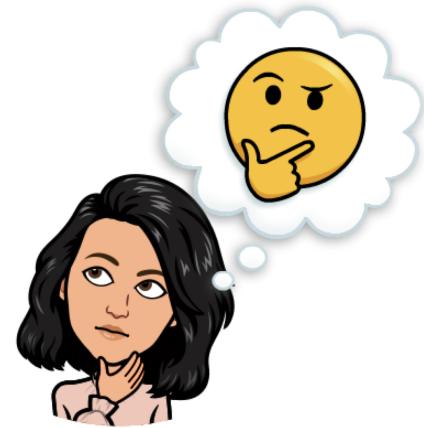
- Now the learner is presented with new images (1 to 4) for prediction.
- What prediction would you expect for each image?

Training data

	CAT
	CAT
	DOG
	DOG

New examples
for prediction

1		?
2		?
3		?
4		?



- Goal: We want the learner to be able to generalize beyond what it has seen in the training data.
- But these new examples should be representative of the training data. That is they should have the same characteristics as the training data.
- In this example, we would like the learner to be able to predict labels for test examples 1, 2, and 3 accurately. Although 2, 3 don't exactly occur in the training data, they are very much similar to the images in the training data. That said, is it fair to expect the learner to label image 4 correctly?

Training error vs. Generalization error

- Given a model M , in ML, people usually talk about two kinds of errors of M .
 - Error on the training data: $\text{error}_{\text{training}}(M)$
 - Error on the entire distribution D of data: $\text{error}_D(M)$
- We are interested in the error on the entire distribution
 - ... But we do not have access to the entire distribution 😞

Data Splitting [video (<https://youtu.be/h2AEobwcUQw>)]

How to approximate generalization error?

A common way is **data splitting**.

- Keep aside some randomly selected portion from the training data.
- `fit` (train) a model on the training portion only.
- `score` (assess) the trained model on this set aside data to get a sense of how well the model would be able to generalize.
- Pretend that the kept aside data is representative of the real distribution D of data.



```
In [9]: 1 # scikit-learn train_test_split
2 url = "https://scikit-learn.org/stable/modules/generated/sklearn.model_
3 IPython.display.IFrame(url, width=1000, height=800)
```

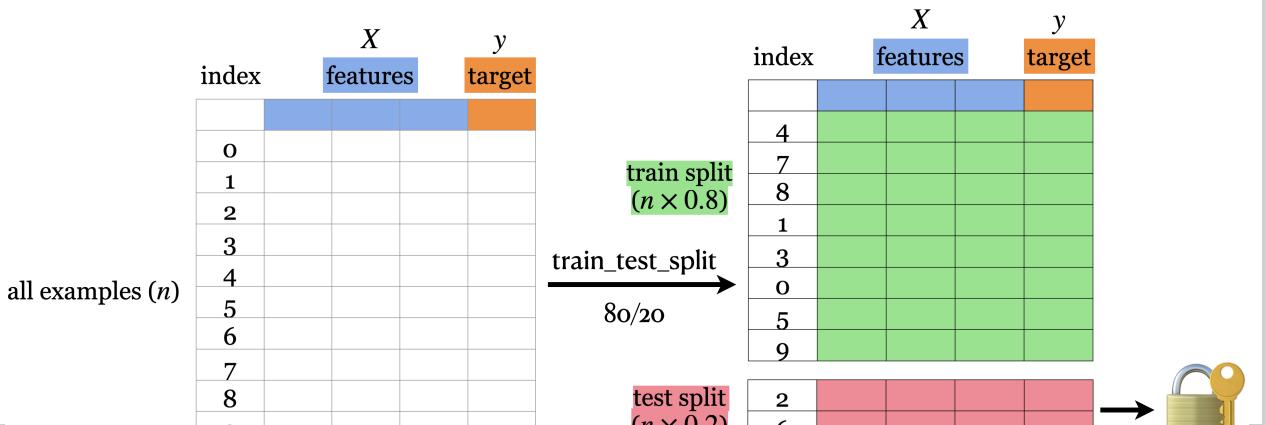
Out[9]:

- We can pass `x` and `y` or a dataframe with both `x` and `y` in it.
- We can also specify the train or test split sizes.

Simple train/test split

- The picture shows an 80%-20% split of a toy dataset with 10 examples.

- The data is shuffled before splitting.
- Usually when we do machine learning we split the data before doing anything and put the test data in an imaginary chest lock.



```
In [10]: 1 # Let's demonstrate this with the canada usa cities data
          2 # The data is available in the data directory
          3 df = pd.read_csv("../data/canada_usa_cities.csv")
          4 X = df.drop(columns=["country"])
          5 y = df["country"]
```

```
In [ ]: 1 X
```

```
In [ ]: 1 y
```

In [11]:

```

1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(
4     X, y, test_size=0.2, random_state=123
5 ) # 80%-20% train test split on X and y
6
7 # Print shapes
8 shape_dict = {
9     "Data portion": ["X", "y", "X_train", "y_train", "X_test", "y_test"]
10    "Shape": [
11        X.shape,
12        y.shape,
13        X_train.shape,
14        y_train.shape,
15        X_test.shape,
16        y_test.shape,
17    ],
18 }
19
20 shape_df = pd.DataFrame(shape_dict)
21 HTML(shape_df.to_html(index=False))

```

Out[11]:

	Data portion	Shape
X	(209, 2)	
y	(209,)	
X_train	(167, 2)	
y_train	(167,)	
X_test	(42, 2)	
y_test	(42,)	

Creating `train_df` and `test_df`

- Sometimes we want to keep the target in the train split for EDA or for visualization.

In []:

```

1 train_df, test_df = train_test_split(
2     df, test_size=0.2, random_state=123
3 ) # 80%-20% train test split on df
4 X_train, y_train = train_df.drop(columns=["country"]), train_df["country"]
5 X_test, y_test = test_df.drop(columns=["country"]), test_df["country"]
6 train_df.head()

```

In []:

```

1 mglearn.discrete_scatter(X.iloc[:, 0], X.iloc[:, 1], y, s=12)
2 plt.xlabel("longitude")
3 plt.ylabel("latitude");

```

```
In [12]: 1 model = DecisionTreeClassifier()
          2 model.fit(X_train.to_numpy(), y_train)
          3 display_tree(X_train.columns, model)
```

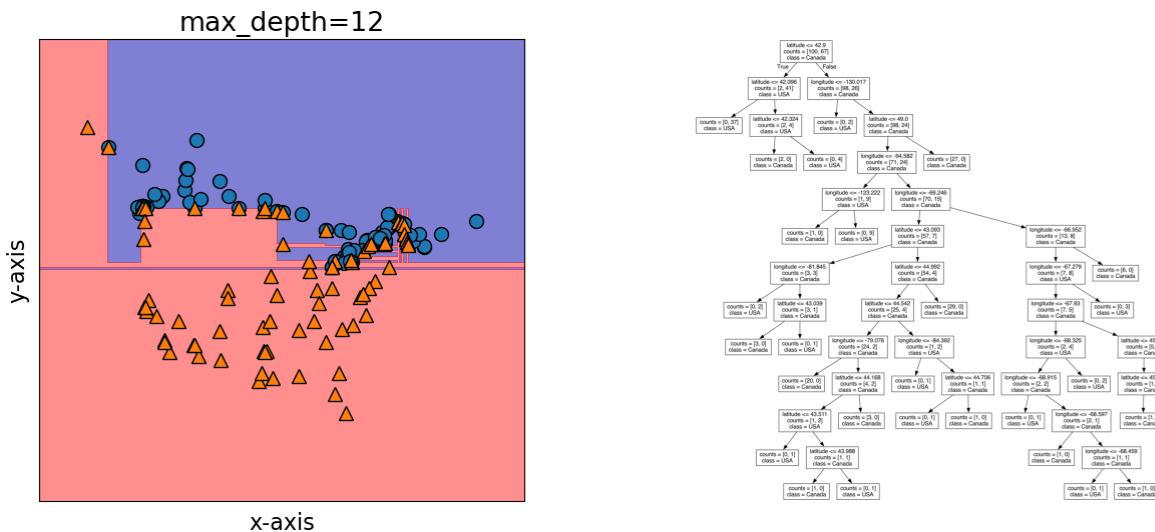
Out[12]: <graphviz.sources.Source at 0x182c5fd90>

Let's examine the train and test accuracies with the split now.

```
In [13]: 1 print("Train accuracy:  %0.3f" % model.score(X_train.to_numpy(), y_train))
          2 print("Test accuracy:   %0.3f" % model.score(X_test.to_numpy(), y_test))
```

Train accuracy: 1.000
Test accuracy: 0.738

```
In [14]: 1 plot_tree_decision_boundary_and_tree(model, x, y, height=6, width=16, e
```

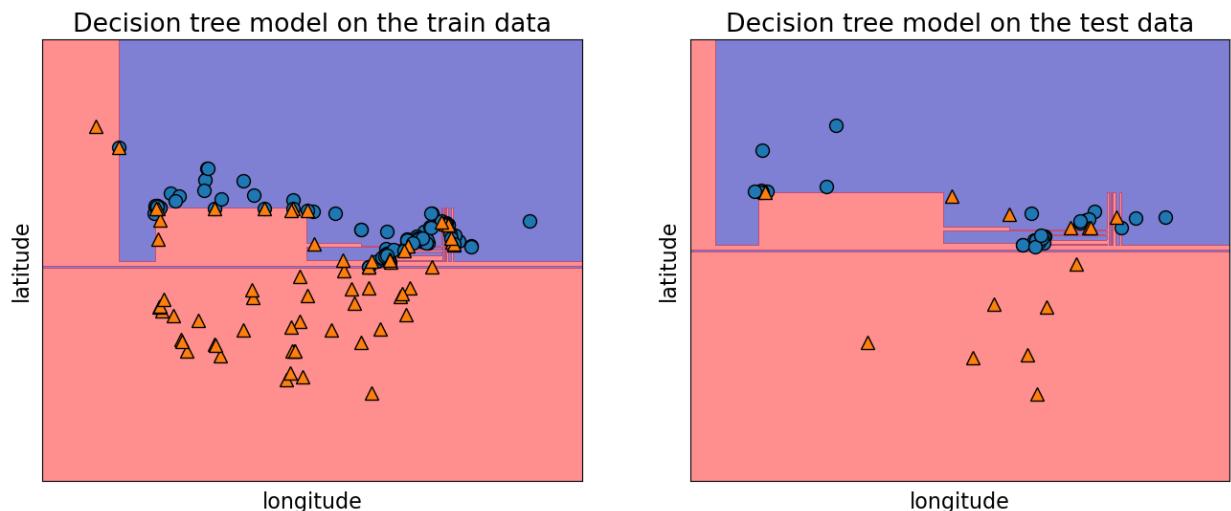


In [15]:

```

1 fig, ax = plt.subplots(1, 2, figsize=(16, 6), subplot_kw={"xticks": ()},
2 plot_tree_decision_boundary(
3     model,
4     X_train,
5     y_train,
6     eps=10,
7     x_label="longitude",
8     y_label="latitude",
9     ax=ax[0],
10    title="Decision tree model on the train data",
11 )
12 plot_tree_decision_boundary(
13     model,
14     X_test,
15     y_test,
16     eps=10,
17     x_label="longitude",
18     y_label="latitude",
19     ax=ax[1],
20    title="Decision tree model on the test data",
21 )

```



What could we do to make the model generalize better to unseen data?

- Useful arguments of `train_test_split`:
 - `test_size`
 - `train_size`
 - `random_state`

`test_size`, `train_size` arguments

- Let's us specify how we want to split the data.
- We can specify either of the two. See the documentation [here](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html) (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html).

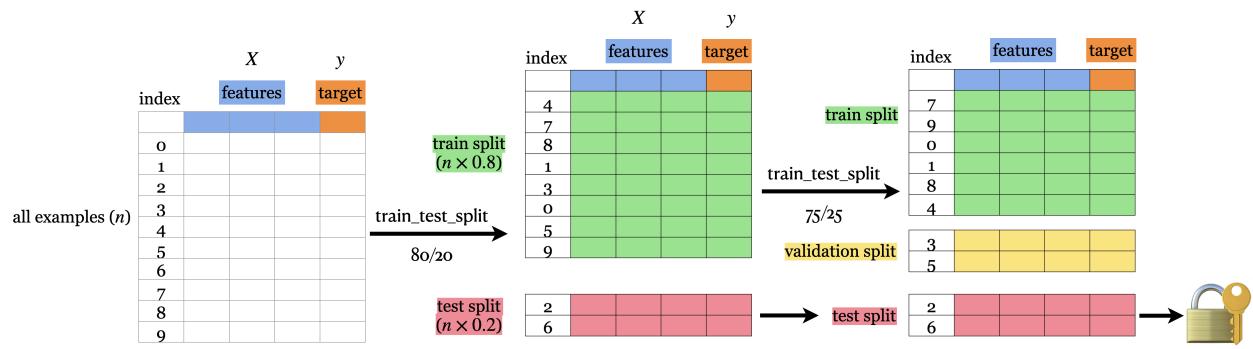
- There is no hard and fast rule on what split sizes should we use.
 - It depends upon how much data is available to you.
- Some common splits are 90/10, 80/20, 70/30 (training/test).
- In the above example, we used 80/20 split.

random_state argument

- The data is shuffled before splitting which is crucial step. (You will explore this in the lab.)
- The `random_state` argument controls this shuffling.
- In the example above we used `random_state=123`. If you run this notebook with the same `random_state` it should give you exactly the same split.
 - Useful when you want reproducible results.

Train/validation/test split

- Some of you may have heard of "validation" data.
- Sometimes it's a good idea to have a separate data for hyperparameter tuning.



- We will try to use "validation" to refer to data where we have access to the target values.
 - But, unlike the training data, we only use this for hyperparameter tuning and model assessment; we don't pass these into `fit`.
- We will try to use "test" to refer to data where we have access to the target values
 - But, unlike training and validation data, we neither use it in training nor hyperparameter optimization.
 - We only use it **once** to evaluate the performance of the best performing model on the validation set.
 - We lock it in a "vault" until we're ready to evaluate.

Note that there isn't good consensus on the terminology of what is validation and what is test.

Validation data is also referred to as **development data** or **dev set** for short.

"Deployment" data

- After we build and finalize a model, we deploy it, and then the model deals with the data in the wild.
- We will use "deployment" to refer to this data, where we do **not** have access to the target values.
- Deployment error is what we *really* care about.
- We use validation and test errors as proxies for deployment error, and we hope they are similar.
- So, if our model does well on the validation and test data, we hope it will do well on deployment data.

Summary of train, validation, test, and deployment data

	fit	score	predict
Train	✓	✓	✓
Validation		✓	✓
Test		once	once
Deployment			✓

You can typically expect $E_{train} < E_{validation} < E_{test} < E_{deployment}$.

?? Questions on generalization and data splitting

1. A decision tree model with no depth is likely to perform very well on the deployment data.
2. Data splitting helps us generalize our model better.
3. Deployment data is used at the very end and only scored once.
4. Validation data could be used for hyperparameter optimization.

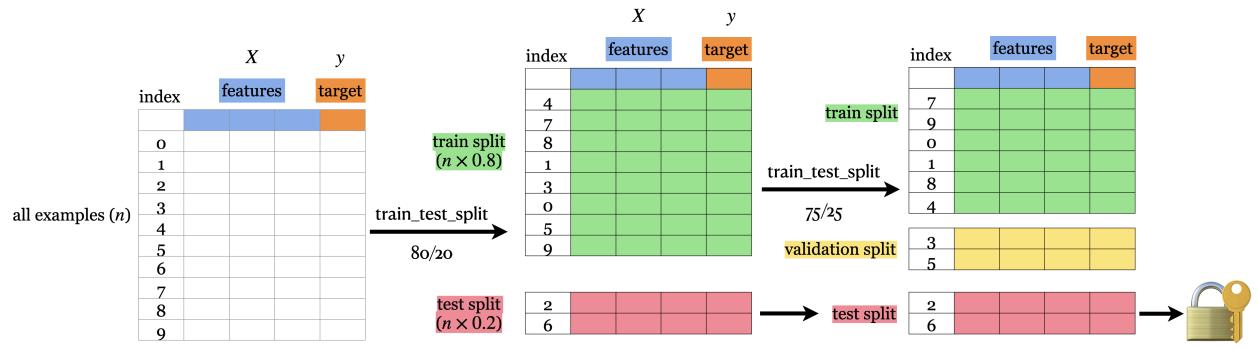
1. False
2. False. Data splitting helps us **assess** how well our model would generalize.
3. False. You cannot score on the deployment data as there are no labels available.
4. True

1. Why you can typically expect $E_{train} < E_{validation} < E_{test} < E_{deployment}$.
2. Discuss the consequences of not shuffling before splitting the data in `train_test_split`.

Cross-validation [video (<https://youtu.be/4cv8VYonepA>)]

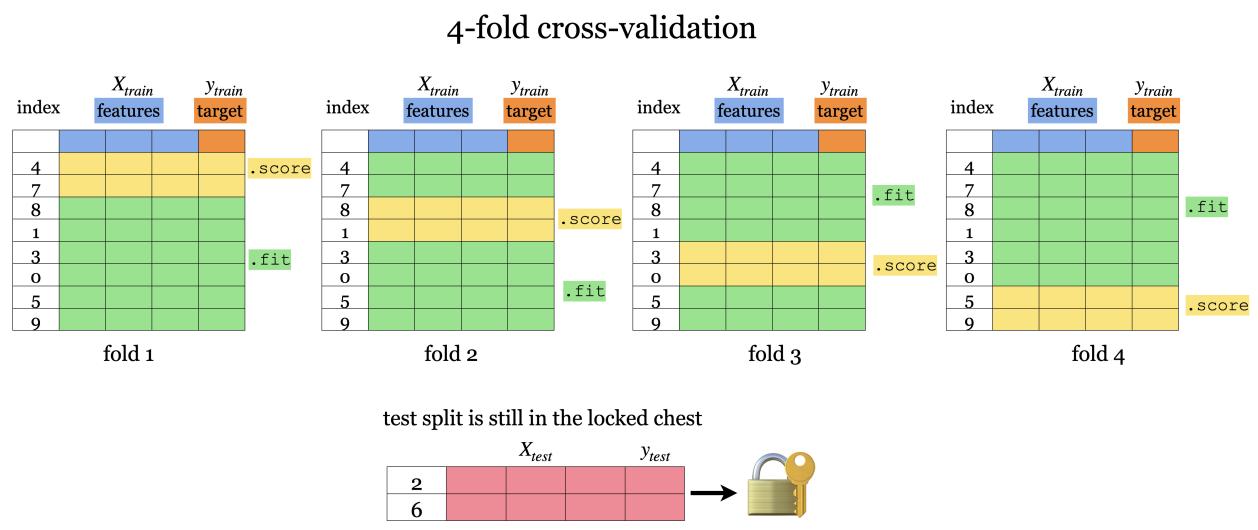
Problems with single train/validation split

- Only using a portion of your data for training and only a portion for validation.
- If your dataset is small you might end up with a tiny training and/or validation set.
- You might be unlucky with your splits such that they don't align well or don't well represent your test data.



Cross-validation to the rescue!!

- Cross-validation provides a solution to this problem.
- Split the data into k folds ($k > 2$, often $k = 10$). In the picture below $k = 4$.
- Each "fold" gets a turn at being the validation set.
- Note that cross-validation doesn't shuffle the data; it's done in `train_test_split`.



- Each fold gives a score and we usually average our k results.
- It's better to examine the variation in the scores across folds.
- Gives a more **robust** measure of error on unseen data.

Cross-validation using `scikit-learn`

```
In [16]: 1 # Let's demonstrate this with the canada usa cities data
          2 # The data is available in the data directory
          3 df = pd.read_csv("../data/canada_usa_cities.csv")
          4 X = df.drop(columns=["country"])
          5 y = df["country"]
```

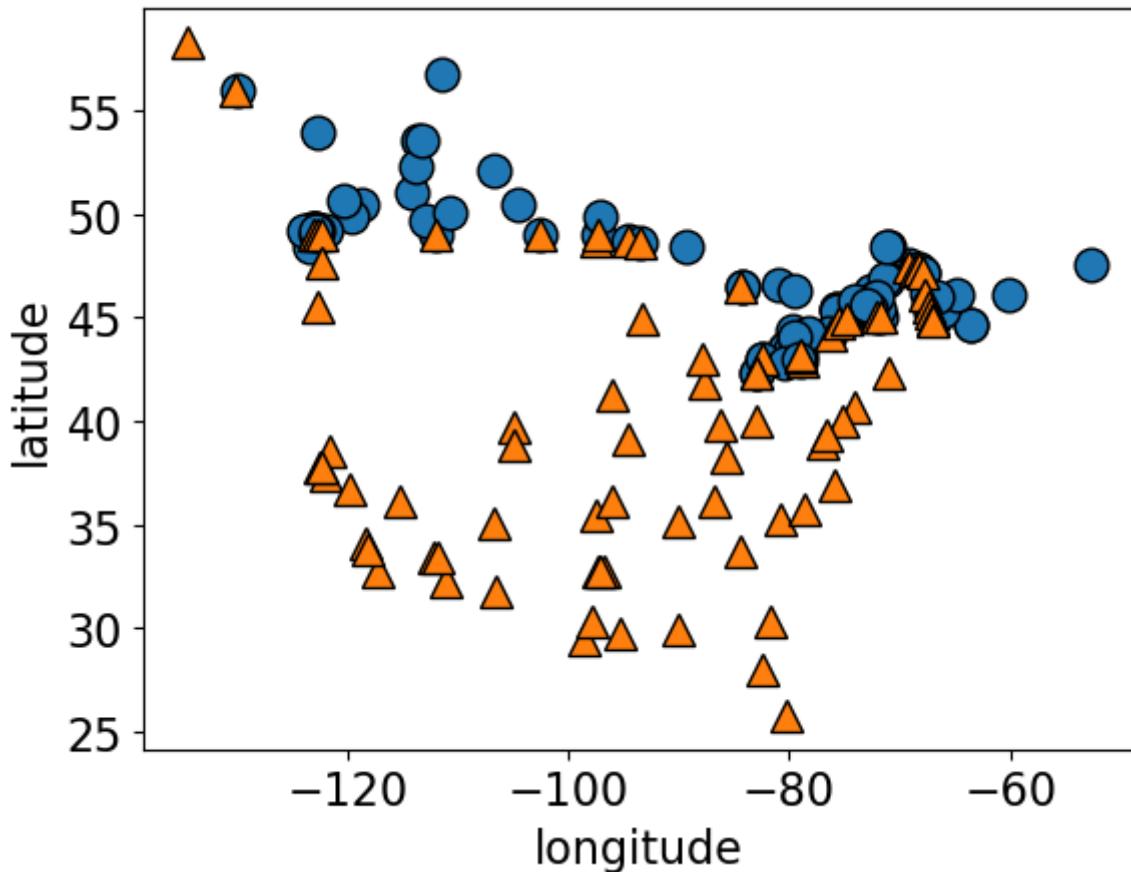
```
In [17]: 1 X
```

Out[17]:

	longitude	latitude
0	-130.0437	55.9773
1	-134.4197	58.3019
2	-123.0780	48.9854
3	-122.7436	48.9881
4	-122.2691	48.9951
...
204	-72.7218	45.3990
205	-66.6458	45.9664
206	-79.2506	42.9931
207	-72.9406	45.6275
208	-79.4608	46.3092

209 rows × 2 columns

```
In [18]: 1 mglearn.discrete_scatter(X.iloc[:, 0], X.iloc[:, 1], y, s=12)
2 plt.xlabel("longitude")
3 plt.ylabel("latitude");
```



```
In [19]: 1 train_df, test_df = train_test_split(
2     df, test_size=0.2, random_state=123
3 ) # 80%-20% train test split on df
4 X_train, y_train = train_df.drop(columns=["country"]), train_df["country"]
5 X_test, y_test = test_df.drop(columns=["country"]), test_df["country"]
```

`cross_val_score`

```
In [20]: 1 from sklearn.model_selection import cross_val_score, cross_validate
```

```
In [21]: 1 model = DecisionTreeClassifier(max_depth=4)
2 cv_scores = cross_val_score(model, X_train, y_train, cv=10)
3 cv_scores
```

```
Out[21]: array([0.76470588, 0.82352941, 0.70588235, 0.94117647, 0.82352941,
 0.82352941, 0.70588235, 0.9375 , 0.9375 , 0.9375 ])
```

```
In [22]: 1 print(f"Average cross-validation score = {np.mean(cv_scores):.2f}")
2 print(f"Standard deviation of cross-validation score = {np.std(cv_scores)}
```

```
Average cross-validation score = 0.84
Standard deviation of cross-validation score = 0.09
```

Under the hood

- It creates `cv` folds on the data.
- In each fold, it fits the model on the training portion and scores on the validation portion.
- The output is a list of validation scores in each fold.

`cross_validate`

- Similar to `cross_val_score` but more powerful.
- Lets us access training and validation scores.

```
In [23]: 1 scores = cross_validate(model, X_train, y_train, cv=10, return_train_sc
2 pd.DataFrame(scores)
```

Out[23]:

	fit_time	score_time	test_score	train_score
0	0.002467	0.001434	0.764706	0.913333
1	0.001965	0.001553	0.823529	0.906667
2	0.001762	0.001211	0.705882	0.906667
3	0.001707	0.001211	0.941176	0.900000
4	0.002209	0.001237	0.823529	0.906667
5	0.001707	0.001201	0.823529	0.913333
6	0.001752	0.001191	0.705882	0.920000
7	0.001724	0.001166	0.937500	0.900662
8	0.001722	0.001112	0.937500	0.900662
9	0.001599	0.001074	0.937500	0.900662

```
In [24]: 1 pd.DataFrame(pd.DataFrame(scores).mean())
```

Out[24]:

	0
fit_time	0.001861
score_time	0.001239
test_score	0.840074
train_score	0.906865

Keep in mind that cross-validation does not return a model. It is not a way to build a model that can be applied to new data. The purpose of cross-validation is to **evaluate** how well the model will generalize to unseen data.

Note that both `cross_val_score` and `cross_validate` functions do not shuffle the data. Check out [StratifiedKFold](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html#sklearn.model_selection.StratifiedKFold) (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html#sklearn.model_selection.StratifiedKFold) where proportions of classes is the same in each fold as they are in the whole dataset. By default, `sklearn` uses `StratifiedKFold` when carrying out cross-validation for classification problems.

In []: 1 `mglearn.plots.plot_cross_validation()`

Our typical supervised learning set up is as follows:

- We are given training data with features `x` and target `y`
- We split the data into train and test portions: `x_train`, `y_train`, `x_test`, `y_test`
- We carry out hyperparameter optimization using cross-validation on the train portion: `x_train` and `y_train`.
- We assess our best performing model on the test portion: `x_test` and `y_test`.
- What we care about is the **test error**, which tells us how well our model can be generalized.
- If this test error is "reasonable" we deploy the model which will be used on new unseen examples.

In [27]: 1 `X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=`
2 `model = DecisionTreeClassifier(max_depth=10)`
3 `scores = cross_validate(model, X_train, y_train, cv=10, return_train_sc`
4 `pd.DataFrame(scores)`
5
6

Out[27]:

	fit_time	score_time	test_score	train_score
0	0.002738	0.001618	0.875000	1.000000
1	0.002283	0.001780	0.875000	0.992857
2	0.001879	0.001215	0.875000	1.000000
3	0.001745	0.001198	0.687500	1.000000
4	0.001719	0.001201	0.812500	1.000000
5	0.002033	0.001563	0.812500	1.000000
6	0.002237	0.001231	0.866667	0.985816
7	0.001750	0.001477	0.600000	0.992908
8	0.002210	0.001361	0.666667	1.000000
9	0.001799	0.001203	0.733333	1.000000

```
In [ ]: 1 def mean_std_cross_val_scores(model, X_train, y_train, **kwargs):
2     """
3         Returns mean and std of cross validation
4     """
5     scores = cross_validate(model, X_train, y_train, **kwargs)
6
7     mean_scores = pd.DataFrame(scores).mean()
8     std_scores = pd.DataFrame(scores).std()
9     out_col = []
10
11    for i in range(len(mean_scores)):
12        out_col.append((f"%0.3f (+/- %0.3f)" % (mean_scores[i], std_scores[i])))
13
14    return pd.Series(data=out_col, index=mean_scores.index)
```

```
In [ ]: 1 results = {}
2 results["Decision tree"] = mean_std_cross_val_scores(
3     model, X_train, y_train, return_train_score=True
4 )
5 pd.DataFrame(results).T
```

? ? Questions on cross-validation

1. k -fold cross-validation calls fit k times. True or False?
2. We use cross-validation to improve model performance. True or False?
3. Discuss advantages and disadvantages of cross-validation.

1. True
2. False. We can use it to assess model performance.

Break (5 min)



Underfitting, overfitting, the fundamental trade-off, the golden rule [[video \(<https://youtu.be/Ihay8yE5KTI>\)](https://youtu.be/Ihay8yE5KTI)]

Types of errors

Imagine that your train and validation errors do not align with each other. How do you diagnose the problem?

We're going to think about 4 types of errors:

- E_{train} is your training error (or mean train error from cross-validation).
- E_{valid} is your validation error (or mean validation error from cross-validation).
- E_{test} is your test error.
- E_{best} is the best possible error you could get for a given problem.

Underfitting

```
In [28]: 1 model = DecisionTreeClassifier(max_depth=1) # decision stump
2 scores = cross_validate(model, X_train, y_train, cv=10, return_train_score=True)
3 print("Train error: %0.3f" % (1 - np.mean(scores["train_score"])))
4 print("Validation error: %0.3f" % (1 - np.mean(scores["test_score"])))
```

Train error: 0.188
 Validation error: 0.212

- If your model is too simple, like `DummyClassifier` or `DecisionTreeClassifier` with `max_depth=1`, it's not going to pick up on some random quirks in the data but it won't even capture useful patterns in the training data.
- The model won't be very good in general. Both train and validation errors would be high. This is **underfitting**.
- The gap between train and validation error is going to be lower.
- $E_{\text{best}} < E_{\text{train}} \lesssim E_{\text{valid}}$

Overfitting

```
In [29]: 1 model = DecisionTreeClassifier(max_depth=None)
2 scores = cross_validate(model, X_train, y_train, cv=10, return_train_score=True)
3 print("Train error: %0.3f" % (1 - np.mean(scores["train_score"])))
4 print("Validation error: %0.3f" % (1 - np.mean(scores["test_score"])))
```

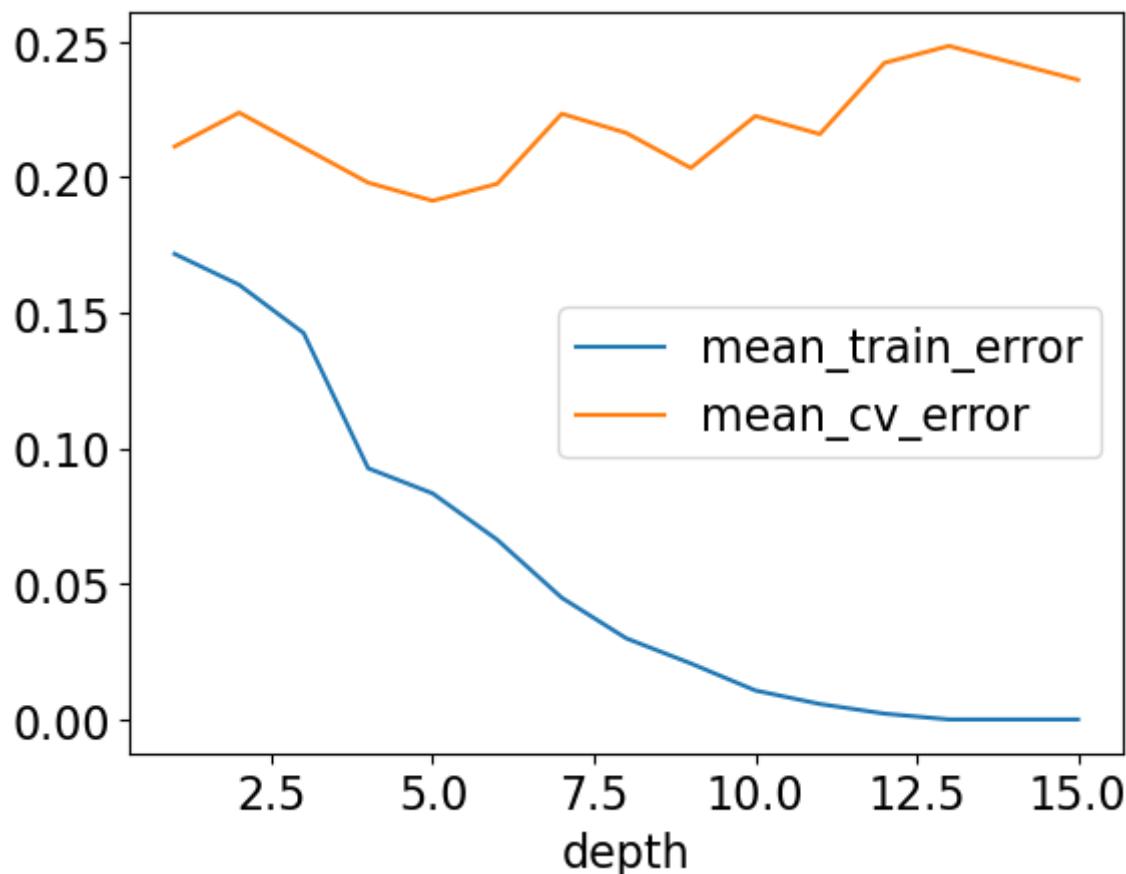
Train error: 0.000
 Validation error: 0.220

- If your model is very complex, like a `DecisionTreeClassifier(max_depth=None)`, then you will learn unreliable patterns in order to get every single training example correct.
- The training error is going to be very low but there will be a big gap between the training error and the validation error. This is **overfitting**.
- In overfitting scenario, usually we'll see: $E_{\text{train}} < E_{\text{best}} < E_{\text{valid}}$
- In general, if E_{train} is low, we are likely to be in the overfitting scenario. It is fairly common to have at least a bit of this.

- So the validation error does not necessarily decrease with the training error.

```
In [30]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
2 results_dict = {
3     "depth": [],
4     "mean_train_error": [],
5     "mean_cv_error": [],
6     "std_cv_error": [],
7     "std_train_error": []
8 }
9 param_grid = {"max_depth": np.arange(1, 16)}
10
11 for depth in param_grid["max_depth"]:
12     model = DecisionTreeClassifier(max_depth=depth)
13     scores = cross_validate(model, X_train, y_train, cv=10, return_train_time=True)
14     results_dict["depth"].append(depth)
15     results_dict["mean_cv_error"].append(1 - np.mean(scores["test_score"]))
16     results_dict["mean_train_error"].append(1 - np.mean(scores["train_score"]))
17     results_dict["std_cv_error"].append(scores["test_score"].std())
18     results_dict["std_train_error"].append(scores["train_score"].std())
19
20 results_df = pd.DataFrame(results_dict)
21 results_df = results_df.set_index("depth")
```

```
In [31]: 1 results_df[["mean_train_error", "mean_cv_error"]].plot();
```



- Here, for larger depths we observe that the training error is close to 0 but validation error goes up and down.
- As we make more complex models we start encoding random quirks in the data, which are not grounded in reality.

- These random quirks do not generalize well to new data.
- This problem of failing to be able to generalize to the validation data or test data is called **overfitting**.

The "fundamental tradeoff" of supervised learning:

As you increase model complexity, E_{train} tends to go down but $E_{\text{valid}} - E_{\text{train}}$ tends to go up.

Bias vs variance tradeoff

- The fundamental trade-off is also called the bias/variance tradeoff in supervised machine learning.

Bias : the tendency to consistently learn the same wrong thing (high bias corresponds to underfitting), or failing to learn something important

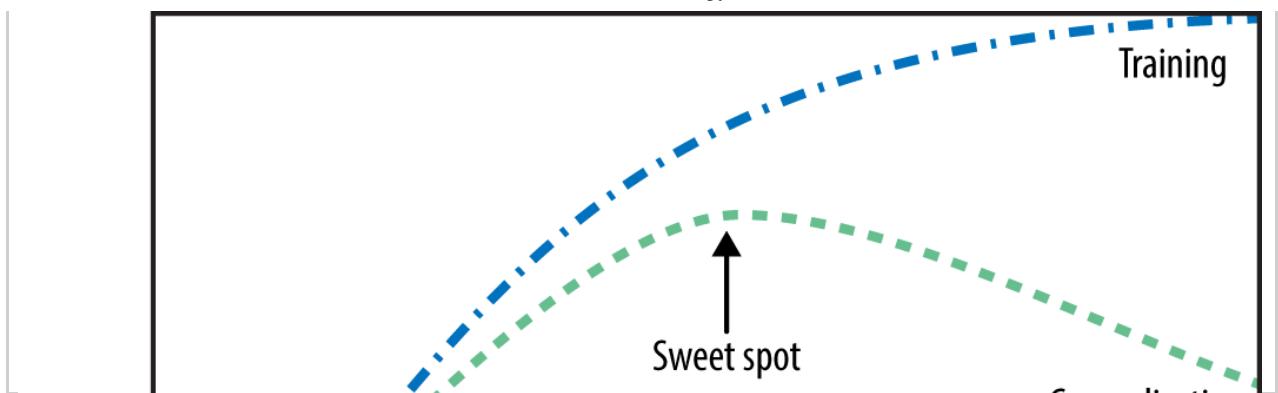
Variance : the tendency to learn random things irrespective of the real signal (high variance corresponds to overfitting)

Check out [this article by Pedro Domingos](#)

(<https://homes.cs.washington.edu/~pedrod/papers/cacm12.pdf>) for some approachable explanation on machine learning fundamentals and bias-variance tradeoff.

How to pick a model that would generalize better?

- We want to avoid both underfitting and overfitting.
- We want to be consistent with the training data but we don't rely too much on it.



- There are many subtleties here and there is no perfect answer but a common practice is to pick the model with minimum cross-validation error.

```
In [ ]: 1 def cross_validate_std(*args, **kwargs):
2     """Like cross_validate, except also gives the standard deviation of
3     res = pd.DataFrame(cross_validate(*args, **kwargs))
4     res_mean = res.mean()
5     res_mean["std_test_score"] = res["test_score"].std()
6     if "train_score" in res:
7         res_mean["std_train_score"] = res["train_score"].std()
8     return res_mean
```

```
In [32]: 1 results_df
```

```
Out[32]:    mean_train_error  mean_cv_error  std_cv_error  std_train_error
```

depth	mean_train_error	mean_cv_error	std_cv_error	std_train_error
1	0.171657	0.211250	0.048378	0.006805
2	0.160258	0.223750	0.062723	0.007316
3	0.142467	0.210833	0.067757	0.022848
4	0.092604	0.197917	0.056955	0.006531
5	0.083338	0.191250	0.067120	0.010650
6	0.066251	0.197500	0.074773	0.012019
7	0.044873	0.223333	0.080734	0.009059
8	0.029909	0.216250	0.088397	0.009422
9	0.020653	0.203333	0.090978	0.010294
10	0.010679	0.222500	0.092669	0.007938
11	0.005699	0.215833	0.109335	0.004264
12	0.002143	0.242083	0.095497	0.003273
13	0.000000	0.248333	0.089485	0.000000
14	0.000000	0.242083	0.086932	0.000000
15	0.000000	0.235833	0.083836	0.000000

test score vs. cross-validation score

```
In [33]: 1 best_depth = results_df.index.values[np.argmin(results_df["mean_cv_error"])
2 print(
3     "The minimum validation error is %0.3f at max_depth = %d "
4     %
5         np.min(results_df["mean_cv_error"]),
6         best_depth,
7     )
8 )
```

The minimum validation error is 0.191 at max_depth = 5

- Let's pick `max_depth = 5` and try this model on the test set.

```
In [34]: 1 model = DecisionTreeClassifier(max_depth=best_depth)
2 model.fit(X_train, y_train)
3 print(f"Error on test set: {1 - model.score(X_test, y_test):.2f}")
```

Error on test set: 0.19

- The test error is comparable with the cross-validation error.
- Do we feel confident that this model would give similar performance when deployed?

The golden rule

- Even though we care the most about test error **THE TEST DATA CANNOT INFLUENCE THE TRAINING PHASE IN ANY WAY.**
- We have to be very careful not to violate it while developing our ML pipeline.
- Even experts end up breaking it sometimes which leads to misleading results and lack of generalization on the real data.

Golden rule violation: Example 1

Emergent Tech ▶ Artificial Intelligence

Was this quake AI a little too artificial? Nature-published research accused of boosting accuracy by mixing training, testing data

Academics, journal deny making a boo boo

By [Katyanna Quach](#) 3 Jul 2019 at 09:01

21 SHARE ▾



Golden rule violation: Example 2



Intelligent Machines

How can we avoid violating golden rule?

- Recall that when we split data, we put our test set in an imaginary vault.



Here is the workflow we'll generally follow.

- Splitting:** Before doing anything, split the data x and y into x_{train} , x_{test} , y_{train} , y_{test} or train_df and test_df using `train_test_split`.
- Select the best model using cross-validation:** Use `cross_validate` with `return_train_score = True` so that we can get access to training scores in each fold. (If we want to plot train vs validation error plots, for instance.)
- Scoring on test data:** Finally score on the test data with the chosen hyperparameters to examine the generalization performance.

Again, there are many subtleties here we'll discuss the golden rule multiple times throughout the course and in the program.

?? Questions for you

Underfitting or overfitting?

1. If the mean train accuracy is much higher than the mean cross-validation accuracy.
2. If the mean train accuracy and the mean cross-validation accuracy are both low and relatively similar in value.
3. Decision tree with no limit on the depth.
4. Decision stump on a complicated classification problem.

1. Overfitting
2. Underfitting
3. Overfitting
4. Underfitting

State whether True/False.

1. In supervised learning, the training error is always lower than the validation error.
2. The fundamental tradeoff of ML states that as training error goes down, validation error goes up.
3. More "complicated" models are more likely to overfit than "simple" ones.
4. If our training error is extremely low, we are likely to be overfitting.
5. A very simple model (e.g. decision stump) has high variance.

1. False
2. False
3. True
4. True
5. False

What did we learn today?

- Importance of generalization in supervised machine learning
- Data splitting as a way to approximate generalization error
- Train, test, validation, deployment data
- Cross-validation
- A typical sequence of steps to train supervised machine learning models
 - training the model on the train split
 - tuning hyperparameters using the validation split
 - checking the generalization performance on the test split

- Overfitting, underfitting, the fundamental tradeoff, and the golden rule.



In []:

1